

(51) Int.Cl. ⁷	識別記号	F I	テマコード* (参考)
H 0 4 L 12/56		G 0 6 F 17/30	4 1 9 A
G 0 6 F 17/30	4 1 9	H 0 4 L 11/20	1 0 2 D

審査請求 未請求 予備審査請求 有 (全 46 頁)

(21) 出願番号 特願2000-512323(P2000-512323)
 (86) (22) 出願日 平成10年5月11日(1998.5.11)
 (85) 翻訳文提出日 平成12年3月15日(2000.3.15)
 (86) 国際出願番号 P C T / S E 9 8 / 0 0 8 5 4
 (87) 国際公開番号 W O 9 9 / 1 4 9 0 6
 (87) 国際公開日 平成11年3月25日(1999.3.25)
 (31) 優先権主張番号 9 7 0 3 3 3 2 - 8
 (32) 優先日 平成9年9月15日(1997.9.15)
 (33) 優先権主張国 スウェーデン (S E)

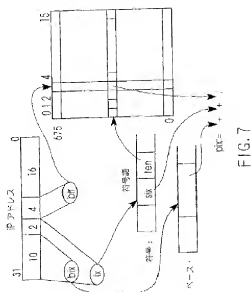
(71) 出願人 エフィシャント ネットワーキング アク
 ティエボラーグ
 スウェーデン国, エス-977 75 リュレ
 オ, アウロルム 2
 (72) 発明者 ブロドレック, アンドレイ
 スロベニア共和国, リュブルヤナ
 1000, ゴツカ 6
 (72) 発明者 デゲルマーク, ミカエル
 スウェーデン国, エス-977 52 リュレ
 オ, アシスレントベーゲン 322
 (74) 代理人 弁理士 石田 敬 (外 4 名)

最終頁に続く

(54) 【発明の名称】 高速ルーティング・ルックアップのための方法とシステム

(57) 【要約】

IPデータグラムをどこに転送するか決定するための、ネクストホップ・テーブル中の関連ネクストホップ情報を有する任意長プレフィックスのエントリを含むルーティング・テーブル中のIPルーティング・ルックアップの方法において、ルーティング・テーブルの表示が、全てのルーティング・テーブル・エントリのプレフィックスによって規定される完全プレフィックス木(7)の形態で保存される。さらに、現在の深さ(D)でのプレフィックス木(7)のカットのデータを含むビット・ベクトル(8)の表示と、ネクストホップ・テーブル及びネクスト・レベル・チャックへの索引を含むポインタの配列とが保存される。前記ビット・ベクトル(8)がビット・マスクに分割され、ビット・マスクの表示がマップテーブルに保存される。その後、各行索引をマップテーブルとポインタ・オフセットに符号化する符号語の配列と、ベース・アドレスの配列とが保存される。最後に、ルックアップが行われる。



【特許請求の範囲】

【請求項1】 IPデータグラムをどこに転送するか決定するための、ネクストホップ・テーブル中の関連ネクストホップ情報を有する任意長プレフィックスのエントリを含むルーティング・テーブル中のIPルーティング・ルックアップの方法であって、

各ノードが子有さないかまたは2つの子を有するかの何れかであり、追加された全ての子が、ネクストホップ情報を有する最も近い先祖と同じネクストホップ情報か、またはこうした先祖が存在しない場合規定されないネクストホップを有する葉であるように完成される、全てのルーティング・テーブル・エントリのプレフィックスによって規定される、完全プレフィックス木(7)の形態での前記ルーティング・テーブルの表示を記憶手段に保存するステップと、

現在の深さ(D)の可能なノード毎に1ビットを有する前記現在の深さの前記プレフィックス木(7)のカットのデータを含むビット・ベクトル(8)の表示を前記記憶手段に保存するステップであって、その際前記プレフィックス木(7)中にノードが存在する場合前記ビットが設定されるステップと、

ポインタの配列と、純粋ヘッドの場合前記ネクストホップ・テーブルへの索引と、根ヘッドの場合ネクスト・レベル・チャンクへの索引とを前記記憶手段に保存するステップと、

前記ビット・ベクトル(8)をある長さのビット・マスクに分割するステップと、

マップテーブル中の可能な前記ビット・マスクの表示を前記記憶手段に保存するステップと、

各々行索引を前記マップテーブルとポインタ・オフセットに符号化する符号語の配列を前記記憶手段に保存するステップと、

ベース・アドレスの配列を前記記憶手段に保存するステップと、

符号語の前記配列中の前記IPアドレスの第1索引部分(ix)に対応する位置の符号語にアクセスするステップと、

前記IPアドレスの列索引部分(bit)と前記マップテーブル中の前記符号語の行索引部分(ten)とに対応する位置のマップテーブル・エントリ部分にアクセスす

るステップと、

ベース・アドレスの前記配列中の前記IPアドレスの第2索引部分(bix)に対応する位置のベース・アドレスにアクセスするステップと、

前記ベース・アドレス・プラス前記符号語のポインタ・オフセット(six) プラス・ポインタの前記配列中の前記マップテーブル・エントリ部分に対応する位置のポインタにアクセスするステップとを特徴とする方法。

【請求項2】 IPデータグラムをどこに転送するか決定するためのネクストホップ・テーブル中の関連ネクストホップ情報を有する任意長プレフィックスのエントリを含むルーティング・テーブル中のIPルーティング・ルックアップのためのシステムであって、各ノードが子を有さないかまたは2つの子を有するかの何れかであり、追加された全ての子が、ネクストホップ情報を有する最も近い先祖と同じネクストホップ情報か、またはこうした先祖が存在しない場合規定されないネクストホップを有する葉である、全てのルーティング・テーブル・エントリのプレフィックスによって規定される、完全プレフィックス木(7)の形態でのルーティング・テーブルと、

現在の深さ(D)の可能なノード毎に1ビットを有する前記現在の深さの前記プレフィックス木(7)のカットのデータを含むビット・ベクトル(8)の表示であって、その際前記プレフィックス木(7)中にノードが存在する場合前記ビットが設定される表示と、

ポインタの配列と、純粋ヘッドの場合前記ネクストホップ・テーブルへの索引と、根ヘッドの場合ネクスト・レベル・チャンクへの索引と、

ある長さのビット・マスクに分割された前記ビット・ベクトル(8)と、

可能な前記ビット・マスクの表示を含むマップテーブルと、

各々行索引を前記マップテーブルとポインタ・オフセットに符号化する符号語の配列と、

ベース・アドレスの配列とを特徴とするシステム。

【発明の詳細な説明】

【0001】

発明の分野

本発明は一般に、IPデータグラムをどこに転送すべきかを決定する、ネクストホップ・テーブル中の関連ネクストホップ情報を有する任意長プレフィックスのエントリを含む、ルーティング・テーブル中のIPルーティング・ルックアップの方法とシステムに関する。

【0002】

従来技術の説明

インターネットはネットワークの相互接続された集合であり、そこでは構成要素であるネットワークは各々そのアイデンティティを保持し、多数のネットワーク相互の通信のためには専用の機構が必要である。構成要素であるネットワークはサブネットワークと呼ばれる。

【0003】

インターネット中の各サブネットワークはそのサブネットワークに接続された装置間の通信をサポートする。さらに、サブネットワークは、網間接続ユニット(IWU)と呼ばれる装置によって接続される。

個々のIWUは、似ている、あるいはそうでない2つのネットワークを接続するために使用されるルータである。このルータは、各ルータとネットワークの各ホストに存在するインターネット・プロトコル(IP)を利用する。

【0004】

IPは局間のコネクションレスまたはデータグラム・サービスを提供する。

ルーティングは一般に、可能な宛先ネットワークの各々について、IPデータグラムを送信すべき次のルータを示すルーティング・テーブルを各局及びルータに維持することで達成される。

IPルータはルーティング・テーブルでルーティング・ルックアップを行い、どこにIPデータグラムを転送するかを決定する。この操作の結果は宛先方向の経路上のネクストホップである。ルーティング・テーブル中のエントリは概念上、関連ネクストホップ情報を有する任意長プレフィックスである。ルーティング・ル

ックアップは、最長一致プレフィックスを有するルーティング・エントリを見出さなければならない。

【0005】

IPルーティング・ルックアップは本質的に低速で複雑だったので、従来技術のソリューションによる操作は、その使用を回避する技術の普及につながった。IPより下の様々なリンク層スイッチング技術、IP層バイパス法（コンピュータ通信会議(IEEE Infocom)会報、カリフォルニア州サンフランシスコ、1996年3月、ギガビット・ネットワーク・ワークショップ会報、ボストン、1995年4月、及びACM SIGCOMM '95、49～58ページ、マサチューセッツ州ケンブリッジ、1995年8月で開示）及びATMのような仮想回線技術に基づく代替ネットワーク層の開発は、ある程度IPルーティング・ルックアップを回避しようという願望の結果である。

【0006】

IPレベルより下のスイッチング・リンク層とフローまたはタグ・スイッチング・アーキテクチャの使用によって、複雑性と冗長性がネットワークに追加される。

最新のIPルータ設計はキャッシュ技術（caching technique）を使用しており、そこでは最近使用された宛先アドレスのルーティング・エントリがキャッシュに保持される。この技術は、トラフィックに十分な局所性が存在するので、キャッシュ・ヒット率が十分に高く、ルーティング・ルックアップの費用がいくつかのパケットにわたって分担されるということに依存している。このキャッシュ法は従来良好に動作した。しかし、現在のインターネットの急速な発達によって必要なアドレス・キャッシュのサイズが増大するに連れて、ハードウェア・キャッシュは不経済なものになることがある。

【0007】

ルーティング・テーブルの伝統的な実現は、ほぼ30年前に発明されたデータ構造であるバトリシア木(ACMジャーナル、15(4):514～534、1968年10月で開示)の、最長プレフィックス一致のために修正したバージョンを使用している。

例えばNetBSD 1.2の実現におけるような、ルーティング・ルックアップ目的でのパトリシア木の直接的な実現は、葉と内部ノードのために24バイトを使用する。40,000エントリの場合、木構造だけでほぼ2メガバイトであり、完全な平衡木では、ルーティング・エントリを見出すために、15あるいは16のノードを横切らなければならない。

【0008】

場合によっては、最長一致プレフィックス規則のため、適切なルーティング情報を見出すために追加ノードの横断が必要になることがあるが、これは初期探索によって適切な葉が見出されることが保証されていないからである。パトリシア木のサイズを縮小しルックアップ速度を改善できる最適化が存在する。それにもかかわらず、データ構造は大きく、それを探索するためにはあまりにも高価なメモリ参照が必要である。現在のインターネット・ルーティング・テーブルは大きすぎてオンチップ・キャッシュには収まらないし、DRAMのオフチップ・メモリ参照は低速すぎて必要なルーティング速度をサポートすることはできない。

【0009】

完全なルーティング・ルックアップを回避することでIPルーティング性能を改善する初期の作業（コンピュータ通信会議(IEEE Infocom)会報、ルイジアナ州ニューオーリンズ、1988年3月で開示）によって見出されたところでは、小さな宛先アドレス・キャッシュによってルーティング・ルックアップ性能は少なくとも65パーセント改善できる。90パーセントを超えるヒット率を得るために必要なスロットは10未満であった。このような小さな宛先アドレス・キャッシュは、現在のインターネットの大きなトラフィック密度とホスト数には不十分である。

【0010】

ATM(非同期転送モード)は、接続セットアップ中にアドレスをネットワークに伝える信号プロトコルを有することでルーティング・ルックアップの実行を回避している。仮想回線識別子(VCI)によってアクセスされる転送状態がセットアップ中に接続の経路に沿ったスイッチにインストールされる。ATMセルは、転送状態を有するテーブルへの直接索引またはハッシュ関数へのキーとして使用される

VCI を含んでいる。ルーティングの決定はATM の場合簡単である。しかし、パケット・サイズが48バイトより大きい場合、さらに多くのATM ルーティングの決定を行う必要がある。

【0011】

タグ・スイッチング及びフロー・スイッチング（コンピュータ通信会議(IEEE Infocom)会報、カリフォルニア州サンフランシスコ、1996年3月で開示）は、ATM 上で動作することを意図する2つのIPパイプ法である。一般的な考え方は、実際のデータ転送を行うリンクレベルATM ハードウェアをIPに制御させるというものである。どのATM 仮想回線識別子を使用し、どのパケットがどのVCI を使用するかをルータ間で一致させる専用プロトコル（コメント要求RFC 1953、インターネット・エンジニアリング・タスク・フォース、1996年5月で開示）が必要である。

【0012】

同じくIP処理の回避を目標とするもう1つのアプローチはIP/ATM アーキテクチャ（ギガビット・ネットワーク・ワークショップ会報、ボストン、1995年4月、及びACM SIGCOMM '95会報、49～58ページ、マサチューセッツ州ケンブリッジ、1995年8月で開示）でなされるが、そこではATM バックプレーンが多数のライン・カードとルーティング・カードを接続する。ルーティング・カードに配置されたIP処理要素がIPヘッダを処理する。パケット・ストリームが到着すると、最初のIPヘッダだけが検討され、その後のパケットは最初のものと同様にルーティングされる。このショートカットの主要な目的は、多数のパケットにわたるIP処理の費用を分担することであると思われる。

【0013】

IPルータ設計はIBM ルータ（高速ネットワークジャーナル、1（4）：281～288、1993年で開示）の場合のように専用ハードウェアを使用してIP処理を行うこともある。これは柔軟性のないソリューションとなる。IPフォーマットまたはプロトコルに何らかの変化があるとこの設計は無効となる。ソフトウェアの柔軟性と汎用プロセッサの急速な性能の向上によってこのソリューションは最適なものになる。ハードウェアによるもう1つのアプローチは、CAM を使用し

てルーティング・ルックアップを行うことである（コンピュータ通信会議(IEEE Infocom)会報、第3巻、1382～1391ページ、サンフランシスコ、1993年で開示）。これは高速ではあるが高価なソリューションである。

【0014】

BBN は現在、転送エンジンとして汎用プロセッサを使用する1対のマルチギガビット・ルータを製造中である。今までのところ発表された情報はほとんどない。しかし、この計画は、転送エンジンとしてアルファ・プロセッサ(Alpha processor)を使用し、全てのIP処理をソフトウェアで行うものと思われる。出版物ギガビット・ネットワーク、マサチューセッツ州レディング、アディソン・ウェズリー社、1993年が示すところによれば、ルート・キャッシュでのヒットを想定すれば、わずか200の命令でIP処理を行うことが可能である。アルファの2次キャッシュは宛先アドレスの大規模LRU(最低使用頻度)キャッシュとして使用される。この方式はトラフィック・パターンの局所性を想定している。局所性が低い場合キャッシュ・ヒット率が低くなりすぎ、性能が犠牲になることがある。

【0015】

発明の概要

従って、ギガビット速度までの各IPパケットについて完全ルーティング・ルックアップを行う改善されたIPルーティング・ルックアップの方法とシステムを提供することが本発明の目的であるが、この方法とシステムは上記で言及された欠点を克服するものである。

【0016】

さらに別の目的は、従来のマイクロプロセッサによるルーティング・ルックアップ速度を向上させることである。

もう1つの目的は、転送テーブルにおけるルックアップ時間を最小化することである。

本発明のもう1つのさらに別の目的は、従来のマイクロプロセッサのキャッシュに完全に収まるデータ構造を提供することである。

【0017】

その結果、メモリ・アクセスは、データ構造が、例えば比較的低速なDRAMからなるメモリに存在する必要がある場合より何倍も高速になる。

これらの目的は本発明によるIPルーティング・ルックアップの方法とシステムによって得られるが、このシステムは、非常に小型の形態で大規模なルーティング・テーブルを表すことができ、わずかなメモリ参照を使用して高速で探索できるデータ構造である。

【0018】

本発明をより詳細に説明し、本発明の利点と特徴を説明するため、好適実施態様が以下詳細に説明され、添付の図面が参照される。

発明の実施の態様

図1を参照すると、ルータ設計は、多数のネットワーク・インバウンド・インタフェース1、ネットワーク・アウトバウンド・インタフェース2、転送エンジン3、及びネットワーク・プロセッサ4を備えているが、これらは全て接続ファブリック5によって相互接続されている。インバウンド・インタフェース1は接続ファブリック5を通じてパケット・ヘッダを転送エンジン3に送信する。一方転送エンジン3はパケットをどの出力インタフェース2に送信すべきかを決定する。この情報はインバウンド・インタフェース1に返送され、そこからパケットがアウトバウンド・インタフェース2に転送される。転送エンジン3の唯一のタスクはパケット・ヘッダを処理することである。ルーティング・プロトコルへの関与、リソースの確保、特別な注意を要するパケットの処理、及び他の管理動作といった他の全てのタスクはネットワーク・プロセッサ4によって処理される。

【0019】

各転送エンジン3は、ネットワーク・プロセッサ4からダウンロードされ、転送エンジン3中の記憶手段に保存されたルーティング・テーブルのローカル・バージョンである転送テーブルを使用してルーティングの決定を行う。ルーティング・アップデートの度に新しい転送テーブルをダウンロードする必要はない。ルーティング・アップデートは頻繁であるが、ルーティング・プロトコルは収束するため若干の時間を要するので、転送テーブルはそれほど陳腐化せず、せいぜい

1秒に1回程度しか変更する必要はない(スタンフォード大学高速ルーティング及びスイッチングに関するワークショップ、1996年12月、http://tiny-tera.stanford.edu/Workshop_Dec96/で開示)。

【0020】

ネットワーク・プロセッサ4は、転送テーブルの高速アップデートと高速生成のために設計された動的ルーティング・テーブルを必要とする。他方、転送テーブルはルックアップ速度について最適化することができ、動的である必要はない。

ルックアップ時間を最小化するために、ルックアップの期間中必要なメモリ・アクセスの数と、データ構造のサイズという2つのパラメータを転送テーブルのデータ構造において最小化しなければならない。

【0021】

メモリ・アクセスは比較的低速で普通ルックアップ手順のボトルネックとなるため、ルックアップの期間中必要なメモリ・アクセスの数を減少させることは重要である。データ構造は十分に小さくすることができれば、従来のマイクロプロセッサのキャッシュに完全に収まる。これは、パトリシア木の場合のように、データ構造が比較的低速なDRAMからなるメモリに存在する必要がある場合より何桁も高速になることを意味する。

【0022】

転送テーブルがキャッシュに完全に収まらないばあいでも、テーブルの大部分がキャッシュに存在することができれば有益である。トラフィック・パターンの局所性によってデータ構造の最も頻繁に使用される部分がキャッシュに保持されるので、大部分のルックアップが高速になる。さらに、少量の必要な外部メモリとして高速SRAMを使用することが実行可能になる。SRAMは高価であり、高速であるほどさらに高価になる。費用が一定の場合、SRAMは必要量が少ない方が高速である。

【0023】

第2の設計目標として、費用のかかる命令と面倒なビット抽出操作を回避するため、データ構造は、ルックアップの期間中に必要な命令が少なく、できる限り

エンティティを自然に整列した状態に保持できるものであるべきである。

データ構造に関する定量的設計パラメータを決定するため、以下説明されるように、多数の大規模ルーティング・テーブルがこれまで検討されている。これらのテーブルに存在する別個のルーティング・エントリは40,000とかなり少ない。ネクストホップが同一であれば、残りのルーティング情報も同じであるので、同じネクストホップを指定する全てのルーティング・エントリはルーティング情報を共有できる。ルータのルーティング・テーブル中の別個のネクストホップの数は、1つのホップで到達できる他のルータまたはホストの数によって制限されるので、大規模バックボーン・ルータの場合でもこの数が小さいのは驚くべきことではない。しかし、ルータが例えば大規模ATMネットワークに接続されている場合、ネクストホップの数はもっと多いことがある。

【0024】

この実施態様では、転送テーブル・データ構造は 2^{14} すなわち16Kの個別ネクストホップに対応するよう設計されているが、これは大部分の場合十分である。別個のネクストホップが256より少ない場合、ネクストホップ・テーブルへの索引は1つのバイトに保存できるので、ここで説明される転送テーブルは、別の実施態様で占有する空間がかなり少なくなるように修正できる。

【0025】

転送テーブルは本質的に3つのレベルを有する木である。1つのレベルの探索には1～4のメモリ・アクセスが必要である。従って、メモリ・アクセスの最大数は12である。しかし、従来のルーティング・テーブルではルックアップの大多数が必要とするのは1～2レベルだけなので、メモリ・アクセス数の大部分は8以下である。

【0026】

データ構造を理解する目的で、図2に示される、IPアドレス空間全体に広がる二分木6を想像のこと。その深さは32であり、葉の数は 2^{32} であるが、これは可能な各IPアドレスについて1つである。ルーティング・テーブル・エントリのプレフィックスはあるノードを末端とする木の中の経路を規定する。そのノードに根付いた部分木の全てのIPアドレス（葉）はそのルーティング・エントリによ

ってルーティングされる。この方法で各ルーティング・テーブル・エントリは同一のルーティング情報を有するIPアドレスの範囲を規定する。

【0027】

いくつかのルーティング・エントリが同じIPアドレスを対象にしている場合、最長一致の規則が適用される。それによれば、あるIPアドレスについて、最長一致プレフィックスを有するルーティング・エントリが使用される。この状況は図3に示されている。ルーティング・エントリ e_1 は、範囲 r 中のアドレスについて e_2 によって隠されている。

【0028】

転送テーブルは、全てのルーティング・エントリがまたがる二分木6、プレフィックス木7の表示である。プレフィックス木は完全であること、すなわち木の各ノードは2つの子を有するかまたは子を有さない何れかであることが必要である。1つの子を有するノードは2つの子を有するように拡張されなければならない。この形で追加された子は常に葉であり、それらのネクストホップ情報は、ネクストホップ情報を有する最も近い先祖のネクストホップと同じであるか、またはそのような先祖が存在しない場合「規定されない」ネクストホップである。

【0029】

図4に示されるこの手順は、プレフィックス木7中のノードの数を増加させるが、小さな転送テーブルの構築を可能にする。しかし、転送テーブルを構築するために実際にプレフィックス木を構築する必要はない。プレフィックス木が使用されるのは説明を簡単にするためである。転送テーブルは、全てのルーティング・エントリを1回通過する間に構築できる。

【0030】

ルーティング・エントリの集合は、IPアドレス空間をIPアドレスの集合に分割する。正しいルーティング情報を見出す問題は、インターバル・セット・メンバーシップ問題 (SIAMコンピュータジャーナル、17 (1) : 1093~1102、1988年12月で開示) と同様である。この場合、各インターバルはプレフィックス木中のノードによって規定されるので、転送テーブルを圧縮するために使用できる特性を有する。例えば、IPアドレスの各範囲は2の累乗である長さを

有する。

【0031】

図5に示されるように、データ構造のレベル1は深さ16までのプレフィックス木を対象とし、レベル2は深さ17から24までを対象とし、レベル3は深さ25から32までを対象とする。プレフィックス木の一部がレベル16の下に延びた場合はいつでも、レベル2チャンク(chunk)が木のその部分を記述する。同様に、レベル3のチャンクはプレフィックス木の24より深い部分を記述する。データ構造のレベルを探索した結果は、ネクストホップ・テーブルへの索引かまたは次のレベルのチャンクの配列への索引である。

【0032】

例えば、図5のデータ構造のレベル1には、深さ16でプレフィックス木7のカット(cut)が存在する。カットは深さ16で可能なノード1つ毎に1ビットのビット・ベクトルに保存される。すなわち、 2^{16} ビット=64Kビット=8Kバイトが必要である。IPアドレスの最初の部分に対応するビットを見出すため、アドレスの上部の16ビットがビット・ベクトルへの索引として使用される。

【0033】

深さ16のプレフィックス木にノードが存在する場合、ベクトル中の対応するビットが設定される。また、木が16未満の深さに葉を有する場合、その葉の対象となるインターバルの最下位ビットが設定される。他の全てのビットは0である。すなわち、ビット・ベクトル中のビットは、

プレフィックス木がカットの下まで続くことを表す1、根ヘッド(root head)(図6のビット6、12及び13)、または、

深さ16またはそれ未満の葉を表す1、純粋ヘッド(genuine head)(図6のビット0、4、7、8、14及び15)、または、

この値が16未満の深さの葉の対象となる範囲のメンバ(member)であることを意味する0(図6のビット1、2、3、5、9、10及び11)である。メンバは前記メンバより小さい最大ヘッドと同じネクストホップを有する。

【0034】

純粋ヘッドの場合、ネクストホップ・テーブルへの索引を保存しなければなら

ない。メンバは前記メンバより小さい最大ヘッドと同じ索引を使用する。根ヘッドの場合、対応する部分木を表すレベル2チャンクへの索引を保存しなければならない。ヘッド情報は配列に保存された16ビット・ポインタで符号化される。ポインタの2ビットはそれがどんな種類のポインタかを符号化し、残りの14ビットはネクストホップ・テーブルまたはレベル2チャンクを含む配列の何れかへの索引である。

【0035】

適切なポインタを見出すために、ビット・ベクトルは長さ16のビット・マスクに分割されるが、それらは $2^{12}=4096$ 存在する。さらに、配列でのポインタの位置は、ベース索引、6ビット・オフセット及び4ビット・オフセットという3つのエンティティを加算することで得られる。ベース索引プラス6ビット・オフセットによって特定のビット・マスクに対応するポインタがどこに保存されるかが決定される。4ビット・オフセットはポインタの中のどれを検索すべきかを指定する。図7は、これらのエンティティを見出す方法を示す。以下の段落はその手順の詳細説明である。

【0036】

ビット・マスクは完全なブレイクツ木から生成されるので、16ビットの全ての組合せが可能なのではない。長さ2nの0でないビット・マスクは、長さnの2つのビット・マスクか、または値1のビット・マスクの何らかの組合せである。a(n)が長さ2ⁿのありうる0でないビット・マスクの数であるとする。a(n)は次の漸化式によって規定される。

【0037】

$$a(0) = 1, a(n) = 1 + a(n-1)^2$$

すなわち、長さ16のありうるビット・マスクの数は $a(4) + 1 = 678$ であるが、1が加算されるのは、ビット・マスクが0のことあるからである。従って各ビット・マスクに対するエントリを有するテーブルへの索引が必要とするのは10ビットだけである。

【0038】

このテーブル、すなわちマップテーブルは、ビット・マスク範囲内のビット数

を4ビット・オフセットにマップするために保持される。このオフセットは必要なポインタを見出すためにいくつポインタをスキップするかを指定するので、ビット索引より小さいセット・ビットの数に等しい。これらのオフセットは、ポインタがたまたまどんな値を有しているかということは無関係に全ての転送テーブルについて同じである。マップテーブルは一定で、一度に全てについて生成される。

【0039】

あらゆるビットマスクは、長さが2の偶数乗であり、同じ2の累乗の倍数であるビット索引で始まるビットのインターバルが1) 全ての0を含むか、または2) 最下位ビット・セットを有するかの何れかであるという特性を有する。

マップテーブルによって提供されるオフセットへのビットマスクとビットからのマッピングは、プレフィックス木を完全にシマッピングの使用を可能にするプレフィックスの拡張と共に、小さい転送テーブルサイズを達成する鍵である。

【0040】

実際のビット・マスクは必要ではなく、ビット・ベクトルの代わりに、マップテーブルへの10ビットの索引プラス6ビット・オフセットからなる16ビット符号語(cord word)の配列を保持する。6ビット・オフセットは64までのポインタを対象とするので、4つの符号語毎に1つのベース索引が必要である。最大64Kのポインタが存在しうるので、ベース索引は最大16ビット($2^{16}=64K$)である必要がある。

【0041】

図7に示される手順では、データ構造の第1レベルを探索するために疑似符号の次のステップが必要であるが、そこでは符号語の配列が符号と呼ばれ、ベース・アドレスの配列がベースと呼ばれ、ixは符号語の配列中のIPアドレスの第1索引部分であり、bitはマップテーブルへのIPアドレスの列索引部分であり、tenはマップテーブルへの符号語の行索引部分であり、bixはベース・アドレスの配列へのIPアドレスの第2索引部分であり、pixは、レベル2チャンクへの索引と共にネクストホップ・テーブルへの索引を含む、ポインタの配列へのポインタである。

【0042】

ix : =IPアドレスの上位12ビット
bit : =IPアドレス符号語の上位16ビットの低位4 : =符号 [ix]
ten : =符号語sixからの10ビット : =符号語からの6ビット
bix : =IPアドレスの上位10ビット
pix : =ベース [bix]+six +マップテーブル [ten] [bit]ポインタ : =レベル1__ポインタ [pix]

すなわち、数ビットの抽出、配列参照及び加算だけが必要である。配列を索引する際の暗黙の乗算以外、乗算または除算命令は必要ない。

【0043】

第1レベルを探索するためにアクセスする必要があるのは合計7バイト、すなわち、2バイトの符号語、2バイトのベース・アドレス、マップテーブル中の1バイト（実際には4ビット）、最後に2バイトのポインタである。第1レベルのサイズは符号語配列について8Kバイト、ベース索引の配列について2Kバイト、プラス多数のポインタである。マップテーブルが必要とする5.3Kバイトは3つのレベル全てで共有される。

【0044】

ビットマスクが0であるかまたは1つのビット集合(bit set)を有する場合、ポインタはネクストホップ・テーブルへの索引でなければならない。こうしたポインタは直接符号語に符号化できるので、マップテーブルはビット・マスク1及び0のためのエントリを含む必要はない。すなわち、マップテーブル・エントリ数は676（索引0～675）まで減少する。符号語（上記のten）の10ビットが675より大きい場合、符号語はネクストホップ・テーブルへの直接索引を表す。符号語からの6ビットが索引の最下位6ビットとして使用され、(ten-676)が索引の上位ビットである。この符号化によって最大(1024-676)×2⁶=22272の索引が可能になるが、これは設計の対象である16Kより多い。この最適化によって、ルーティング・エントリが深さ12以上にある時3つのメモリ参照が除去され、ポインタ配列のポインタの数がかなり減少する。これは比較及び条件付き分岐による。

【0045】

マッピングはCプログラミング言語の以下のデータと関数によって示される。
この符号の重要な特徴は、それが提供する個別のビットマスクからオフセットの配列へのマッピングである。

mtable (及びmt) のエントリを入れ替えるといった変化があっても与えられるマッピングは同じである。従って、4つの16ビット語の代わりに2つの32ビット語といった、オフセットの配列を表す他の方法もあるだろう。

/* マップテーブル */

#define MAPVECLEN 4

typedef uint16 MAPVEC [MAPVECLEN] ;

typedef MAPVEC MCOMPACT;

MCOMPACT mt [TMAX];

/* mtはルックアップ時に使用されるマップテーブルである。

【0046】

構築中に使用されるmtableから初期化される。 */

```
typedef struct mentry {
    uint16 mask; /* 16ビット・パターン (LSB は
                  設定!) */
    uint16 len; /* 8ビットlen (値1~16) */
    MAPVEC map; /* 4のグループ中の16の4ビットオフセッ
                  ト */
}
```

} MENTRY, *MP;

void mentry2mcompact(MENTRY *from, MCOMPACT *to) /* 「from」からマップ部分を取り、それを「to」に置く */

```
{
    int i;
    for(i=0; i<MAPVECLEN; i++) {
        (*to)[i]=from->map[i];
    }
}
```

```

}
extern void mtable__compact()
/* mtableからmtを初期化する */
{
    register int i;
    for(i=0; i<TMAX; i++) {
        mentry2mcompact(&mtable[i], &mt[i]);
    }
}

```

データ構造のレベル2及び3はチャンクからなる。チャンクは高さ8の部分木を対象とし、最大 $2^8 = 256$ ヘッドを含みうる。レベル $n-1$ の根ヘッドはレベル n のチャンクを指す。

【0047】

イマジナリ・ビット・ベクトルに含まれるヘッドの数によって、チャンクには3つの種類がある。すなわち、1～8のヘッドが存在する場合、チャンクは疎であり、ヘッドの8ビット索引の配列、それに加えて8つの16ビット・ポインタ、すなわち合計24バイトによって表される。

9～64のヘッドが存在する場合、チャンクは稠密である。それは、ベース索引の数以外レベル1と同様に表される。相違点は、6ビット・オフセットが64ポインタ全てに及ぶため、16ビット符号語全てについて必要なベース索引は1つだけだということである。合計34バイト、それに加えてポインタのための18～128バイトが必要である。

【0048】

65～256のヘッドが存在する場合、チャンクは超稠密である。それはレベル1と同様に表される。16の符号語と4のベース索引で合計40バイトとなる。さらに65～256のポインタが130～512バイトを必要とする。

稠密及び超稠密チャンクは第1レベルと同様に探索される。

疎なチャンクは、8要素用によく合わせた専用二分探索によって探索される。これは線形探索及び汎用二分探索よりかなり高速である。さらに、この探索は、

条件付き移動命令を有するプロセッサ・アーキテクチャ上で従来の条件付き飛び越しを使用せずに実現できる。

```
/* 疎なチャンクのための探索関数 */
static inline uint16 findsparse(SPARSECHUNK *chu, uint32 val)
/* chu->vals は8ビット値のソートされた配列0..7である。
```

【0049】

```
chu->rinfoはポインタの対応する配列0..7である。
val は探索キーである
```

```
*/
{
    uint8 *p, *q;
    p=q&;(chu->vals [0 ] );
    p+=( *(p-3)>val)<<2;
    p+=( *(p+1)>val)<<1;
    p+=( *p>val);

    return (chu->rinfo [p-q ] );
}
```

稠密及び超稠密チャンクは、説明されたようにレベル1と同様に最適化される。疎なチャンクでは、2つの連続ヘッ드의ネクストホップが同一な場合それらを併合 (merge) し、小さい方で表すことができる。チャンクが疎か稠密かを決定する場合、この併合が考慮されるので、併合されたヘッドの数が8以下の場合チャンクは疎であると考えられる。木を完全にするために追加された葉の多くは順番に発生し、同一のネクストホップを有する。こうした葉に対応するヘッドは併合されて疎なチャンクになる。

【0050】

この最適化によって、チャンクの分布はさらに大きい稠密なチャンクからさらに小さい疎なチャンクの方向にシフトされる。大きなテーブルの場合、転送テーブルのサイズは通常5～15パーセント縮小される。

このデータ構造はルーティング・エントリがかなり増大しても対応できる。現在の設計には2つの制限がある。

【0051】

1. 各種類のチャンクの数にレベル当たり 2^{14} 16384 に制限される。

表1は、これが現在使用されているものより約16倍大きいことを示している。それでもこの制限を越える場合、ポインタの符号化を変更し索引にさらに余裕を与えるようにするか、またはポインタ・サイズを増大するようにデータ構造を修正できる。

【0052】

2. レベル2及び3のポインタの数はベース索引のサイズによって制限される。

現在の実現は16ビット・ベース索引を使用し、3～5の増大係数に対応できる。この限度を越える場合、ベース・ポインタのサイズを3バイトに増大するのが簡単である。チャンク・サイズは稠密チャンクの場合3パーセント、超稠密チャンクの場合10パーセント増大する。疎なチャンクには影響はない。

【0053】

このデータ構造が、ルーティング・エントリの数の大きな増加に対応できることは明らかである。そのサイズはルーティング・エントリの数と共にほぼ直線的に増大する。

転送テーブルの性能を調べるために、多数のIPルーティング・テーブルが集められた。インターネット・ルーティング・テーブルは現在インターネット・パフォーマンス・メジャーメント・アンド・アナリシス (IPMA) プロジェクトのウェブサイト (<http://www.ra.net/statistics/>) で入手可能であり、以前はルーティング・アービタ・プロジェクト (<http://www.ra.net/statistics/>) で入手可能となっていたが、これは現在終了している。収集されたルーティング・テーブルは、様々な大規模インターネット相互接続点で使用するルーティング・テーブルの日々のスナップショットである。これらのテーブルのルーティング・エントリの中には多数のネクストホップを含むものもある。この場合、それらの1つが、転送テーブルで使用するネクストホップとして無作為に選択された。

【0054】

図8の表1は様々なルーティング・テーブルから構成された転送テーブルに関するデータを示す。各サイトについて、この表は最大転送テーブルを生成したルーティング・テーブルに関するデータと結果を示す。ルーティング・エントリはルーティング・テーブル中のルーティング・エントリの数であり、ネクストホップはテーブル中に見られる個別のネクストホップの数である。葉はプレフィックス木を完全にするために葉が追加された後のプレフィックス木中の葉の数である。

【0055】

表1の構築時間は、ルーティング・テーブルのメモリ内二分木表示から転送テーブルを生成するために必要な時間を示す。時間はDEC OSF1を実行する333MHzアルファ21164で測定された。次の列は生成されたテーブルでの疎、稠密及び超稠密チャンクの合計数を示し、それに、データ構造の最下位レベルのチャンクの数が続く。

【0056】

表1から、新しい転送テーブルが急速に生成できることが明らかである。1Hzの再生周波数では、消費されるアルファの能力は10分の1未満である。上記で説明されたように、1Hzより高い再生周波数は必要ない。

表1のさらに大きなテーブルはアルファの96Kバイト2次キャッシュには完全に収まらない。しかし、2次キャッシュに収まらない部分のための第3レベル・キャッシュに少量の超高速SRAMを有し、2次キャッシュのミスの費用を低減することは実行可能である。トラフィック・パターンの局所性によって、大部分のメモリ参照は2次キャッシュへのものとなる。

【0057】

観察された興味深い点は、これらのテーブルのサイズが、配列中の全てのプレフィックスをちょうど保存するために要するものと同等だということである。さらに大きなテーブルでは、プレフィックス毎に必要なのはわずか5.6バイトに過ぎない。これらのバイトの半分以上はポインタによって消費される。スプリント・テーブル(Sprint table)では、33469のポインタが存在し、65Kバ

ト以上の記憶装置を必要とする。ポインタの数を減らすことで転送テーブルのサイズをさらに縮小できることは明らかである。

【0058】

ルックアップ・ルーチンに関する測定は、GNU C コンパイラgcc(gnu cc. の使用と移植マニュアル、フリー・ソフトウェア・ファウンデーション、1995年11月、ISBN 1-882114-66-3で開示) によってコンパイルされたC関数上でなされる。報告された時間には、ネクストホップ・テーブルへの関数呼び出しまたはメモリ・アクセスは含まれない。gcc は、最悪の場合データ構造の1レベルを探索するために約50のアルファ命令を使用する符号を生成する。ペンティアム・プロでは、gcc は最悪の場合レベル毎に35から45の命令を使用する符号を生成する。

【0059】

以下のC符号関数は、測定で使われるルックアップ関数の符号を示す。

レベル1符号語配列はint1と呼ばれ、ベース索引配列はbase1と呼ばれる。レベル1のポインタは配列htablに保存される。

チャンクは配列cis、cid、ciddに保存され、ここでiはレベルである。

稠密及び超稠密チャンクに関するベース・アドレスとポインタはそれぞれ配列baseid、baseidd、及びhtabid、htabiddに保存される。

/* 転送テーブルのためのルックアップ関数 */

```
#include "conf.h"
```

```
#include "forward.h"
```

```
#include "mtable.h"
```

```
#include "bit2index.h"
```

```
#include "sparse.h"
```

```
#include "timing.h"
```

```
#include "lookup.h"
```

```
/* これらのマクロは、ipが32ビット符号なしint(uint32) である時だけ動作
```

```

する */
#define EXTRACT(start, bits, ip) (((ip)<<(start))>>(32-(bits)))
#define GETTEN(m) (((m)<<22)>>22)
#define GETSIX(m) ((m)>>10)
#define bit2o(ix, bit)
((mt [(ix)] [(bit)>>2] >>)((bit)&0x3)<<2)&0xf)
/* ルックアップ(ipaddr)--ipaddrに関するルーティング・テーブル・エントリ
への索引 */
unsigned int lookup(uint32 ipaddr)
{
uint32 ix;          /* 符号語配列への索引 */
uint32 code;        /* 16ビット符号語 */
int32 diff;         /* TMAXとの差 */
uint32 ten;         /* mtへの索引 */
uint32 six;         /* 符号の6つの「余分の」ビット */
int32 nhop;         /* ネクストホップへのポインタ */
uint32 off;         /* ポインタ・オフセット */
uint32 hbase;       /* ハッシュ・テーブルのためのベース */
uint32 pntr;        /* ハッシュ・テーブル・エントリ */
int32 kind;         /* pntrの種類 */
uint32 chunk;       /* チャンク索引 */
uint8 *p,*q;        /* 探索球へのポインタ */
uint32 key;         /* 球の探索キー */
/* *HERE* からタイミングを取る */

【0060】
【表1】

```

```

ix = EXTRACT(0,12,ipaddr);
code = int1[ix];
ten = GETTEN(code);
six = GETSIX(code);

if ((diff=(ten-TMAX))>=0) {

    nhop = (six | (diff<<6));
    goto end;
}
off = bit2o(ten, EXTRACT(12,4,ipaddr));
hbase = base1[ix>>2];
pntr = htab1[hbase+six+off];

if ((kind = (pntr & 0x3)) {

    chunk = (pntr>>2);

    if (--kind) {

```

【0061】

【表2】


```

if (--kind) {

    key = EXTRACT(16,8,ipaddr);
    p = q = c2s[chunk].vals[0];

    if ( *(p+3) > key) {
        p += 4;
    };

    while( *p > key) {
        p++;
    };

    pntr = c2s[chunk].rinfo[p - q];

} else {

    ix = EXTRACT(16,4,ipaddr);
    code = c2dd[chunk][ix];
    ten = GETTEN(code);
    six = GETSIX(code);

    if ((diff=(ten-TMAX))>=0) {

        nhop = {six | (diff<<6)};
        goto end;
    }
    hbase = htab2ddbase[(chunk<<2)|(ix>>2)];
    off = bit2o(ten, EXTRACT(20,4,ipaddr));
    pntr = htab2dd[hbase+six+off];
}
} else {

    ix = EXTRACT(16,4,ipaddr);
    code = c2d[chunk][ix];

```

【0062】

【表3】

```

ten = GETTEN(code);
six = GETSIX(code);

if ((diff=(ten-TMAX))>=0) {

    nhop = (six | (diff<<6));
    goto end;
}
hbase = htab2dbase(chunk);
off = bit2o(ten, EXTRACT(16,4,ipaddr));
pntr = htab2d(hbase+six+off);
}

if ((kind = (pntr & 0x3)) {

    chunk = (pntr>>2);

    if (--kind) {
        If (--kind) {

            key = EXTRACT(24,8,ipaddr);
            p = q = c3s[chunk].vals[0];

            if ( *(p+3) > key) {
                p += 4;
            };

            while( *p > key) {
                p++;
            };

            pntr = c3s[chunk].rinfo[p - q];

        } else {

```

【0063】

【表4】

```

ix = EXTRACT(24,4,ipaddr);
code = c3dd[chunk][ix];
ten = GETTEN(code);
six = GETSIX(code);

if ((diff=(ten-TMAX))>=0) {

    nhop = (six | (diff<<6));
    goto end;
}

off = bit2o(ten, EXTRACT(28,4,ipaddr));
hbase = htab3ddbbase[(chunk<<2) | (ix>>2)];
pntr = htab3dd[hbase+six+off];
}

} else {

ix = EXTRACT(24,4,ipaddr);
code = c3d[chunk][ix];
ten = GETTEN(code);
six = GETSIX(code);

if ((diff=(ten-TMAX))>=0) {

    nhop = (six | (diff<<6));
    goto end;
}

off = bit2o(ten, EXTRACT(28,4,ipaddr));
hbase = htab3dbbase[chunk];
pntr = htab3d[hbase+six+off];
}

}

nhop = (pntr >> 2);

end:

【0064】
/* HERE! までタイミングを取る */
return nhop;
}

```

アルファ及びペンティアム・プロ(Pentium Pro)のクロックサイクル・カウンタの現在の値を読むことが可能である。高い精度でルックアップ時間を測定するためにこの機構が使用された。1クロックチック(clock tick)は200MHzで

は5ナノ秒であり、333MHzでは3ナノ秒である。

【0065】

理想的には、転送テーブル全体がキャッシュに配置され、ルックアップは非擾乱キャッシュ(undisturbed cache)によって行われる。これは専用転送エンジンのキャッシュの挙動をエミュレートする。しかし、測定は従来の汎用ワークステーションで行われたのであって、こうしたシステムでキャッシュ内容を制御することは困難である。キャッシュは、I/Oが行われる時、割込みが発生する時、または他の処理の実行が開始される時常に擾乱される。キャッシュを擾乱せずには、測定データをプリントアウトすることやファイルから新しいIPアドレスを読み出すことも不可能である。

【0066】

使用される方法は各ルックアップを2回行い、2回目のルックアップのルックアップ時間を測定する。この方法では、最初のルックアップは擾乱キャッシュ(disturbed cache)によってなされ、2回目は全ての必要なデータが最初のルックアップによって1次キャッシュに置かれた状態のキャッシュでなされる。各1組のルックアップの後で測定データがプリントアウトされ新しいアドレスがフェッチされるが、この手順は再びキャッシュを擾乱する。

【0067】

2回目のルックアップは、データと命令がプロセッサに最も近い1次キャッシュに移動しているため、転送エンジンにおけるルックアップより良好に行われる。ルックアップ時間の上限値を得るために、2次キャッシュへのメモリ・アクセスのために必要な追加時間を測定時間に加算しなければならない。転送テーブルを通る全ての経路を試験するために、完全木への拡張によって追加されたエントリを含む、ルーティング・テーブルの各エントリについてルックアップ時間が測定された。

【0068】

現実のトラフィック・ミックス(traffic mix)が各ルーティング・エントリへのアクセスに関して均一な確率を有することはありそうにないので、平均ルックアップ時間はこれらの実験からは推論できない。さらに、トラフィック・パター

ンの局所性によってデータ構造の頻繁にアクセスされる部分は1次キャッシュに保持されるので、平均ルックアップ時間は減少する。以下計算される性能値は、全てのメモリ・アクセスは1次キャッシュでミスし、常に最悪の場合の実行時間が発生することを想定しているため控えめなものである。現実のルックアップ速度はもっと高速である。

【0069】

表1は、データ構造のレベル3ではチャンクが非常に少ないことを示す。このため、ネクストホップを見出すためルックアップの大部分が探索する必要があるのはせいぜい2レベルまでということになる。従って、2次キャッシュへのメモリ・アクセスのための追加時間は最悪の場合の12ではなく8のメモリ・アクセスについて計算される。もしルックアップの大きな部分がそれらの少数のチャンクにアクセスする場合はあれば、それらは1次キャッシュに移動するので、12全てのメモリ・アクセスはより費用のかからないものになるだろう。

【0070】

実験はアルファ21164上で、333MHzのクロック周波数で行われた。1サイクルは3ナノ秒である。8Kバイト1次データ・キャッシュへのアクセスは2サイクルで完了し、2次96Kバイト・キャッシュへのアクセスには8サイクルが必要である。図9の表2を参照のこと。

図10は、1月1日からのスプリント・ルーティング・テーブルに関するアルファの場合の2回目のルックアップの期間に経過したクロック・チックの分布を示す。観察された最高速のルックアップは17クロック・サイクルを必要とする。これは、第1レベルの符号語がネクストホップ・テーブルへの索引を直接符号化する場合である。こうしたルーティング・エントリは非常に少数である。しかし、こうしたルーティング・エントリは各々多くのIPアドレスを対象とするので、実際のトラフィックはこうした宛先アドレスを多く含むことがある。ルックアップには22サイクルかかるものがあるが、これは前と同じ場合であろう。クロック・サイクル・カウンタが2つの連続命令について読み取られる場合、その差は予想される0ではなく5サイクルのことがあることが、実験によって確認されている。

【0071】

図10の次のスパイク(spike)は41クロック・サイクルにあるが、これは第1レベルにあるポインタがネクストホップ・テーブルへの索引である場合である。伝統的なクラスBアドレスがこのカテゴリに入る。52～53、57、62、67、及び72チックのスパイクは、疎なレベル2チャンクで1、2、3、4または5の値を調べた後にポインタが見出されることに対応する。75及び83チックのスパイクが非常に大きいのは、多くのチックがそれぞれ稠密及び超稠密チャンクを探索する必要があるためである。83以上でいくつかのチックが観察されるのは、おそらく実行時間の変化のために疎なレベル3チャンクを探索した後で見出されるポインタに対応する。2次キャッシュでのキャッシュの競合、またはルックアップ前のパイプライン及びキャッシュ・システムの状態の差によってこうした変化が起きることがある。100クロック・サイクルを越える観察値の末尾は、こうした変化かまたはキャッシュ・ミスの何れかによるものである。全てのデータが1次キャッシュにある場合300ナノ秒あれば十分である。

【0072】

1次キャッシュと2次キャッシュのデータ・アクセスの差は $8 - 2 = 6$ サイクルである。データ構造の2レベルを探索するのに必要なクロック・サイクルは最悪の場合、図10に示された場合より $8 \times 6 = 48$ 多い。これは、2レベルが十分である時、最悪の場合のルックアップでは最大 $100 + 48 = 148$ サイクルすなわち444ナノ秒が必要となることを意味する。すなわち、アルファは、2次キャッシュ中の転送テーブルで、1秒当たり220万のルーティング・ルックアップを行うことができる。

【0073】

もう1つの実験はペンティアム・プロ上で、200MHzのクロック周波数で行われた。1サイクルは5ナノ秒である。1次8Kバイト・キャッシュは2サイクルの待ち時間を有し、256Kバイトの2次キャッシュは6サイクルの待ち時間を有する。表2を参照のこと。

図11は、アルファ21164の場合と同じ転送テーブルに関するペンティアム・プロの場合の2回目のルックアップの期間に経過したクロック・チックの分

布を示す。クロック・サイクル・カウンタを取り出す (fetch)一連の命令は33クロック・サイクルを要する。互いの直後に2つの取り出しが発生する場合カウンタ値は33異なる。この理由で、報告された時間は全て33減らされている。

【0074】

観察された最高速のルックアップは11クロック・サイクルであり、アルファ21164の場合とほぼ同じ速度である。ネクストホップ索引が第1レベルの直後にある場合に対応するスパイクは25クロック・サイクルで発生する。疎なレベル2チャンクに対応するスパイクは36~40クロック・サイクルの範囲に互いに接近して集まっている。ペンティアムの異なったキャッシュ構造は、アルファ21164のキャッシュ構造より線形走査を良好に扱うように思われる。

【0075】

第2レベル・チャンクが稠密及び超稠密である時、ルックアップはそれぞれ48及び50サイクルを必要とする。69まではいくつかの付加的な不規則スパイクがあるが、それ以上で観察されるものは非常に少ない。全てのデータが1次キャッシュにある場合、ルックアップを行うには69サイクル(345ナノ秒)で十分なことは明らかである。

【0076】

1次キャッシュと2次キャッシュのアクセス時間の差は20ナノ秒(4サイクル)である。2つのレベルを調べる必要がある時、ペンティアム・プロの場合のルックアップ時間は最悪で $69 + 8 \times 4 = 101$ サイクルすなわち505ナノ秒である。ペンティアム・プロは2次キャッシュの転送テーブルで1秒当たり少なくとも200万のルーティング・ルックアップを行うことができる。

【0077】

本発明が、上記で示した目標と利点を完全に満足する改善されたIPルーティング・ルックアップの方法とシステムを提供することは明らかであろう。本発明は特定の実施態様と共に説明されたが、代替案、修正及び変形は当業技術分野に熟練した者には明らかである。

上記で説明された実施態様ではルックアップ関数はCプログラミング言語で実装されている。プログラミングの技術分野に熟練した者には、他のプログラミング

グ言語も使用できることが明らかであろう。また、ルックアップ関数は標準デジタル設計技術を使用してハードウェアで実現することもできるが、これもハードウェア設計の技術分野に熟練した者には明らかであろう。

【0078】

例えば、本発明は、アルファ 21164 またはペンティアム・プロを使用するシステム以外のコンピュータ・システム構成、プレフィックス木をカットする他の方法、木の様々なレベル数、マップテーブルを表す他の方法、及び符号語を符号化する他の方法に適用可能である。

さらに、本発明はファイアウォール・ルーティング・ルックアップにも利用できる。

【図面の簡単な説明】

【図1】

図1は、ルータ設計の概略図である。

【図2】

図2は、IPアドレス空間全体にわたる二分木を示す図である。

【図3】

図3は、IPアドレスの範囲を規定するルーティング・エントリを示す図である。

【図4】

図4は、プレフィックス木を拡張して完全なものにするステップを示す図である。

【図5】

図5は、本発明によるデータ構造の3つのレベルを示す図である。

【図6】

図6は、深さ16でのプレフィックス木のカットの一部を示す図である。

【図7】

図7は、データ構造の第1レベル探索を示す図である。

【図8】

図8は、様々なルーティング・テーブルから構成された転送テーブル上のデー

タを示す表である。

【図9】

図9は、プロセッサとキャッシュのデータを示す表である。

【図10】

図10は、アルファ21164に関するルックアップ時間分布を示すグラフである。

【図11】

図11は、ペンティアム・プロに関するルックアップ時間分布を示すグラフである。

【図1】

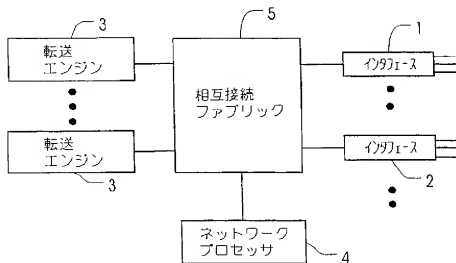


FIG. 1

【図2】



FIG. 2

【図3】

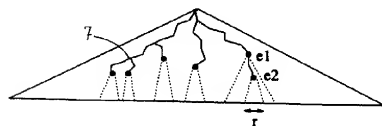


FIG. 3

【図4】

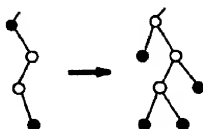


FIG. 4

【図5】

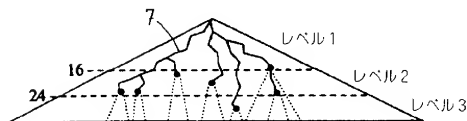


FIG. 5

【图6】

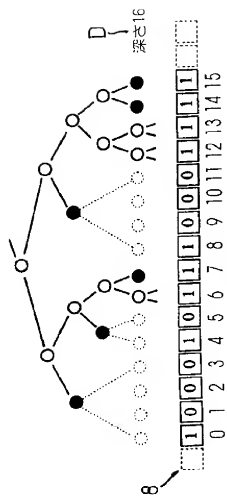


FIG. 6

【図7】

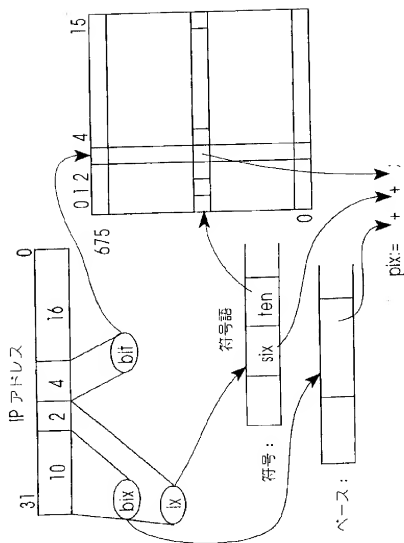


FIG. 7

表1

サイト	月日	年	ルーティング エントリー	ネットワーク ホップ	サイズ (Kb)	構築 時間	疎な チャネル	秘密 チャネル	超秘密 チャネル	レベル3 チャネル
Mae E	1月9日	'97	32732	56	160	99 ミリ秒	1199	587	186	2
Mae E	10月21日	'96	38141	50	148	91 ミリ秒	1060	593	149	4
Sprint	1月1日	'97	21797	17	123	72 ミリ秒	988	483	98	3
Pac	1月28日	'97	18308	2	99	49 ミリ秒	873	357	67	0
Bell										
MaeW	1月1日	'97	12049	28273	51	86	46 ミリ秒	775	312	3
AADS	1月4日	'97	1109	5670	12	28	11 ミリ秒	320	38	0

FIG. 8

表2

プロセッサ	クロック サイクル	1次キャッシュ		2次キャッシュ		3次キャッシュ	
		サイズ	待ち時間	サイズ	待ち時間	サイズ	待ち時間
アルファ21164	3ナノ秒	8 K/バイト	6ナノ秒	96 K/バイト	24ナノ秒	2 M/バイト	2ナノ秒
ハインティアプロ	5ナノ秒	8 K/バイト	10ナノ秒	256 K/バイト	30ナノ秒		

FIG. 9

【図 1 0】

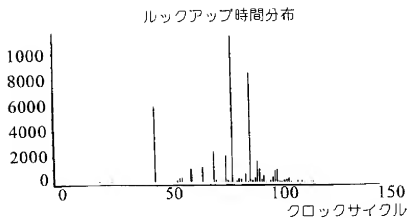


FIG. 10

【図 1 1】

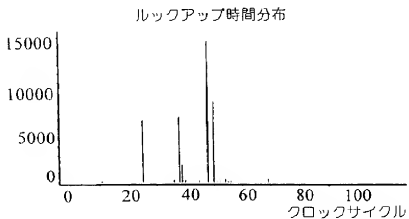


FIG. 11

【手続補正書】特許協力条約第34条補正の翻訳文提出書

【提出日】平成12年3月15日(2000.3.15)

【手続補正1】

【補正対象書類名】明細書

【補正対象項目名】特許請求の範囲

【補正方法】変更

【補正内容】

【特許請求の範囲】

【請求項1】 IPデータグラムをどこに転送するか決定するための、ネクストホップ・テーブル中の関連ネクストホップ情報を有する任意長プレフィックスのエントリを含むルーティング・テーブル中のIPルーティング・ルックアップの方法であって、

各ノードが子を有さないかまたは2つの子を有するかの何れかであり、追加された全ての子が、ネクストホップ情報を有する最も近い先祖と同じネクストホップ情報か、またはそうした先祖が存在しない場合規定されないネクストホップを有する葉であるように完成される、全てのルーティング・テーブル・エントリのプレフィックスによって規定される、完全プレフィックス木(7)の形態での前記ルーティング・テーブルの表示を記憶手段に保存するステップと、

現在の深さ(D)の可能なノード毎に1ビットを有する前記現在の深さの前記プレフィックス木(7)のカットのデータを含むビット・ベクトル(8)の表示を前記記憶手段に保存するステップであって、その際前記プレフィックス木(7)中にノードが存在する場合前記ビットが設定されるステップと、

純粹ヘッドの場合前記ネクストホップ・テーブルへの索引を含み、根ヘッドの場合ネクスト・レベル・チャンクへの索引を含む、ポインタの第1レベル配列を前記記憶手段に保存するステップと、

前記ビット・ベクトル(8)をある長さのビット・マスクに分割するステップと、

マップテーブル中の可能な前記ビット・マスクの表示を前記記憶手段に保存するステップと、

各々行索引を前記マップテーブルとポインタ・オフセットに符号化する符号語の配列を前記記憶手段に保存するステップと、

ベース・アドレスの配列を前記記憶手段に保存するステップと、

符号語の前記配列中の前記IPアドレスの第1索引部分(ix)に対応する位置の符号語にアクセスするステップと、

前記IPアドレスの列索引部分(bit)と、前記マップテーブル中の前記符号語の行索引部分(ten)とに対応する位置のマップテーブル・エントリ部分にアクセスするステップと、

ベース・アドレスの前記配列中の前記IPアドレスの第2索引部分(bix)に対応する位置のベース・アドレスにアクセスするステップと、

前記ベース・アドレス・プラス前記符号語のポインタ・オフセット(six) プラス・ポインタの前記配列中の前記マップテーブル・エントリ部分に対応する位置のポインタにアクセスするステップとを含む方法において、

マップテーブル中の可能な前記ビット・マスクの表示を前記記憶手段に保存する前記ステップが、

長さが2の偶数乗であり、同じ2の累乗の倍数であるビット索引で始まるビットのどのインターバルも全ての0を含むか、または最下位ビット・セットを有するかの何れかであるという特性を有する可能な前記ビット・マスクの表示を保存するステップを含むことを特徴とする方法。

【請求項2】 前記ポインタがネクストレベル・チャンクへの索引である場合、さらに、純粹ヘッドの場合前記ネクストホップ・テーブルへの索引を含み、根ヘッドの場合第3レベル・チャンクへの索引を含む、ポインタの第2レベル・チャンク配列中のポインタを見出すために、第2レベル・チャンクにアクセスするステップを含むことを特徴とする、請求項1に記載の方法。

【請求項3】 ポインタの前記第2レベル・チャンク配列中の前記ポインタが前記第3レベル・チャンクへの索引である場合、さらに、前記ネクストホップ・テーブルへの索引を含むポインタの第3レベル・チャンク配列中のポインタを見出すために、前記第3レベル・チャンクにアクセスするステップを含むことを特徴とする、請求項2に記載の方法。

【請求項4】 前記チャンクが1～8ヘッドを含む場合、前記チャンクが疎なチャンクであり、前記ヘッドの8ビット索引の配列と8つの16ビット・ポイントによって表されることと、

前記チャンクが9～64ヘッドを含む場合、前記チャンクが稠密なチャンクであり、前記第1レベルと同様に表されることと、

前記チャンクが65～256ヘッドを含む場合、前記チャンクが超稠密なチャンクであり、前記第1レベルと同様に表されることを特徴とする、請求項1～3の何れか1項に記載の方法。

【請求項5】 IPデータグラムをどこに転送するか決定するためのネクストホップ・テーブル中の関連ネクストホップ情報を有する任意長プレフィックスのエントリを含むルーティング・テーブル中のIPルーティング・ルックアップのためのシステムであって、各ノードが子を有さないかまたは2つの子を有するかの何れかであり、追加された全ての子、ネクストホップ情報を有する最も近い先祖と同じネクストホップ情報か、またはこうした先祖が存在しない場合規定されないネクストホップを有する葉である、全てのルーティング・テーブル・エントリのプレフィックスによって規定される、完全プレフィックス木(7)の形態でのルーティング・テーブルと、

現在の深さ(D)の可能なノード毎に1ビットを有する前記現在の深さの前記プレフィックス木(7)のカットのデータを含むビット・ベクトル(8)の表示であって、その際前記プレフィックス木(7)中にノードが存在する場合前記ビットが設定される表示と、

純粋ヘッドの場合前記ネクストホップ・テーブルへの索引を含み、根ヘッドの場合ネクストレベル・チャンクへの索引を含む、ポイントの配列と、

ある長さのビット・マスクに分割された前記ビット・ベクトル(8)と、

可能な前記ビット・マスクの表示を含むマップテーブルと、

各々行索引を前記マップテーブルとポイント・オフセットに符号化する符号番の配列と、

ベース・アドレスの配列とを備える方法において、

可能な前記ビット・マスクが、長さが2の偶数乗であり、同じ2の累乗の倍数

であるビット索引で始まるビットのどのインターバルも全ての0を含むか、または最下位ビット・セットを有するかの場合であるという特性を有することを特徴とするシステム。

【請求項6】 チャンクの前記ネクストレベルが、純粹ヘッドの場合前記ネクストホップ・テーブルへの索引を含み、根ヘッドの場合第3レベル・チャンクへの索引を含むポインタの第2レベル・チャンク配列によって表される第2レベル・チャンクであることを特徴とする、請求項5に記載のシステム。

【請求項7】 前記第3レベル・チャンクが、前記ネクストホップ・テーブルへの索引を含む、ポインタの第3レベル・チャンク配列によって表されることを特徴とする、請求項6に記載のシステム。

INTERNATIONAL SEARCH REPORT

International application No.
PCT/SE 98/00854

A. CLASSIFICATION OF SUBJECT MATTER		
IPC6: H04L 12/56 According to International Patent Classification (IPC) or to both national classification and IPC		
B. FIELDS SEARCHED		
Minimum documentation searched (classification system followed by classification symbols)		
IPC6: H04L		
Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched SE,UK,FI,NO classes as above		
Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)		
EDOC, WPIL, JAPIO, INSPEC		
C. DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	Mikael Degermark, "FREE ASPECTS OF PACKET FORWARDING IN THE INTERNET", April 1997, (Luleå Tekniska Högskola, Sweden), page 27 - page 44, Doctoral thesis --	1,2
A	Infocom, Volume 3, 1993, (Morristown, USA) Anthony J. McAuley & Paul Francis, "Fast Routing Table Lookup Using CAMs" --	1,2
A	IBM Technical Disclosure Bulletin, Volume 36, No 2, February 1993, ., "Memory Organization Scheme for the Implementation of Routing Tables in High Performance IP Routers" page 151 - page 153 --	1,2
<input checked="" type="checkbox"/> Further documents are listed in the continuation of Box C. <input type="checkbox"/> See patent family annex.		
<p>* Special categories of cited documents</p> <p>"A" document defining the general state of the art which is not considered to be of particular relevance</p> <p>"B" earlier document but published on or after the international filing date</p> <p>"I" document which may throw doubt on priority claim or which is cited to establish the publication date of another citation or other special reason (as specified)</p> <p>"O" document referring to an oral disclosure, use, exhibition or other means</p> <p>"P" document published prior to the international filing date but later than the priority date claimed</p> <p>"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention</p> <p>"X" document of particular relevance: the claimed invention cannot be considered as not or cannot be considered to involve an inventive step when the document is taken alone</p> <p>"Y" document of particular relevance: the claimed invention cannot be considered to involve an inventive step when the document is considered with one or more other such documents, such combination being obvious to a person skilled in the art</p> <p>"Z" document: member of the same patent family</p>		
Date of the actual completion of the international search		Date of mailing of the international search report
19 October 1998		25 -10- 1998
Name and mailing address of the ISA/ Swedish Patent Office Box 5055, S-102 42 STOCKHOLM Facsimile No. +46 8 666 02 86		Authorized officer Anders Ströbeck Telephone No. +46 8 782 25 00

Form PCT/ISA/210 (second sheet) (July 1994)

INTERNATIONAL SEARCH REPORT

International application No.

PCT/SE 98/00854

C (Continuation). DOCUMENTS CONSIDERED TO BE RELEVANT		
Category*	Quotation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	IBM Disclosure Bulletin, Volume 40, No 3, March 1997, . . , "Technique for Performing Generalized Prefix Matches" page 189 - page 200 --	1,2
A	SIAM Journal Computers, Volume 17, No 6, December 1988, Kurt Mehlhorn et al, "A LOWER BOUND ON THE COMPLEXITY OF THE UNION-SPLIT-FIND PROBLEM" page 1093 - page 1102 -- -----	1,2

Form PCT/ISA/210 (continuation of second sheet) (July 1992)

(81)指定国 EP(AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OA(BF, BJ, CF, CG, CI, CM, GA, GN, ML, MR, NE, SN, TD, TG), AP(GH, GM, KE, LS, MW, SD, SZ, UG, ZW), EA(AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GE, GH, GM, GW, HU, ID, IL, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, UZ, VN, YU, ZW

(72)発明者 カールソン, スバンテ
スウェーデン国, エス-977 53 リュレ
オ, ニュステイーゲン 4

(72)発明者 ビンク, ステファン
スウェーデン国, エス-165 57 ハース
レビー, ビベカ トロレス グランド 8